

ارائه‌ی الگوریتمی به منظور خوشه‌بندی صفحات وب براساس محتوا و لینک

محمد فتحیان* (استاد)

امیرحسین کریمی مجد (دانشجوی دکتری)

دانشکده‌ی مهندسی صنایع، دانشگاه علم و صنعت ایران

مهندسی صنایع و مدیریت شریف، تابستان ۱۳۹۶ (۱۳۹۶)
دوری (۳۳-۱) شماره ۱/۱، ص. ۲۸-۲۱

وجود یک موتور جست‌وجوی کارا می‌تواند سبب افزایش رضایت کاربران از خدمات تحت وب باشد. چالش اصلی موتورهای جست‌وجو، انتخاب مناسب‌ترین صفحات در مواجهه با پرسش‌های چندوجهی کاربران است. «خوشه‌بندی صفحات براساس محتوا و لینک» رویکردی است که برای حل چنین مسائلی در ادبیات پیشنهاد شده است. در این نوشتار، بر یکی از الگوریتم‌های موجود، به نام CohsMix^۱، تمرکز شده و این الگوریتم برای ارتقای کیفیت پاسخ‌ها و افزایش سرعت حل بهبود داده شده است. تعیین نقطه‌ی شروع مناسب، استفاده از خواص شبکه‌های پیچیده به منظور ساده‌سازی محاسبات، و محاسبه‌ی مقدار واقعی انحراف استاندارد از جمله تغییرات پیشنهادی برای بهبود الگوریتم است. نتایج تجربی نشان می‌دهد که الگوریتم بهبودیافته، کیفیت جواب‌ها را ارتقا داده و باعث افزایش سرعت حل می‌شود. همچنین، به‌عنوان مطالعه‌ی موردی، داده‌های مربوط به وبلاگ‌های فارسی استخراج و الگوریتم بهبودیافته روی این داده‌ها اجرا خواهد شد.

واژگان کلیدی: خوشه‌بندی، تجارت الکترونیکی، محتوا، لینک، موتور جست‌وجو، شبکه‌های پیچیده.

۱. مقدمه

در دهه‌ی اخیر شاهد رشد و توسعه‌ی چشمگیری در حوزه‌ی تجارت الکترونیکی بوده‌ایم. در این حوزه، اگر چه فروریختن مرزهای جغرافیایی و دسترسی آسان و حتی رایگان به منابع مختلف موجب اقبال کاربران به خدمات الکترونیکی شده، اما سرعت و دقت در دسترسی به منابع مورد نظر در میان انبوه صفحات وب نزد آن‌ها اهمیت ویژه‌ی دارد. به عبارت دیگر، یکی از عوامل مؤثر برای تأمین رضایت کاربران از خدمات تحت وب، وجود ابزار مناسبی برای یافتن منابع مورد نظر است. بدین ترتیب، برای یک پایگاه تحت وب، برخورداری از یک موتور جست‌وجوی کارا می‌تواند حرکتی در راستای جذب مشتریان و افزایش رضایت‌مندی آنها باشد.

مسئله‌ی که موتورهای جست‌وجو عمدتاً با آن روبرو هستند این است که خروجی آن‌ها در مواجهه با پرسش‌های^۲ چندوجهی کاربران، طیف گسترده و پراکنده‌ی از صفحات خواهد بود و انتخاب مناسب‌ترین و مرتبط‌ترین صفحات کار دشواری است. به عبارت ساده‌تر احتمال این که خروجی موتور جست‌وجو در نظر کاربران نامرتب جلوه کند، بسیار زیاد است. به‌عنوان مثال واژه‌ی «شیر» حداقل در سه معنا کاربرد دارد: ۱. حیوانی وحشی؛ ۲. شیرگاو؛ ۳. شیرآب. حال زمانی که کاربری شیر به معنای شیر آب را مورد جست‌وجو قرار دهد و نتایج جست‌وجو حول دو کاربرد دیگر (که در مقایسه با شیر آب بیشتر مورد استفاده قرار می‌گیرند) ارائه شود، موتور

* نویسنده مسئول

تاریخ: دریافت ۱۳۹۳/۳/۲۷، اصلاحیه ۱۳۹۴/۴/۷، پذیرش ۱۳۹۴/۵/۷.

fathian@iust.ac.ir
karimimajd@iust.ac.ir

جست‌وجو نتوانسته رضایت کاربر را تأمین کند. تقریباً از ابتدای ظهور وب، برای حل چنین معضلی رویکرد «خوشه‌بندی^۳ موضوعی صفحات» پیشنهاد شد. بدین ترتیب برای هر مفهوم متفاوت، یک خوشه و از هر خوشه نیز یک نماینده در نتایج جست‌وجو قابل رؤیت خواهد بود. این رویکرد در ابتدا مورد استقبال واقع شد، اما در ادامه مشخص شد که کارایی لازم را ندارد؛ زیرا قادر نبود گروه‌های متنوعی را که در ارتباط با کاربردهای مختلف واژگان چندبعدی شکل گرفته‌اند از یکدیگر تمیز دهد.

در سال ۲۰۱۰ الگوریتمی به نام CohsMix معرفی شد^۱ که خوشه‌بندی صفحات را براساس محتوا^۴ و لینک^۵ میان آن‌ها انجام می‌داد. برتری این روش نسبت به پژوهش‌های انجام شده تا قبل از سال ۲۰۱۰^{۱،۲} این است که آن‌ها با در نظر گرفتن مدل شکل‌گیری لینک‌های میان صفحات، امکان استخراج الگوهای واقعی میان صفحات را مهیا کردند. همچنین در نظر گرفتن توزیع آماری مناسب برای مقادیر مربوط به ویژگی‌های صفحات، به مدل آن‌ها -- در مقایسه با مدل‌های مبتنی بر فاصله^{۱،۲} -- انعطاف‌پذیری بیشتری می‌دهد. علی‌رغم چنین قابلیت‌های مهمی، حجم بالای محاسبات لازم برای اجرای این الگوریتم به زمان‌گیر بودن آن منجر شده است.

در این مقاله ابتدا بر مشکل نحوه‌ی انجام محاسبات الگوریتم CohsMix و ارائه‌ی الگوریتمی برای بهبود آن تمرکز خواهد شد. ایده‌ی اصلی الگوریتم بهبودیافته در مقایسه با الگوریتم پایه، بهره‌گیری بیشتر از ویژگی‌های ماتریس مجاورت و اطلاعات

مربوط به توزیع احتمالات خصیصه‌ها در راستای کاهش زمان حل مسئله و افزایش کیفیت پاسخ‌هاست. به‌عنوان مطالعه‌ی موردی، برای اولین بار داده‌های وبلاگ‌های فارسی را استخراج کرده و به خوشه‌بندی آن‌ها به‌منظور فراهم آوردن بستری برای ارتقای خروجی موتورهای جست‌وجوگر خواهیم پرداخت.

در بخش بعدی به پیشینه‌ی پژوهش می‌پردازیم و در بخش ۳ مدل پیشنهادی تشریح خواهد شد. بخش ۴ نتایج عددی حاصل از آزمایش‌های مختلف روی مجموعه داده‌های موجود در ادبیات را ارائه می‌دهد. در بخش ۵ نتیجه‌گیری کلی ارائه خواهد شد.

۲. پیشینه‌ی پژوهش

در این بخش ابتدا پژوهش‌های انجام‌شده در حوزه‌ی خوشه‌بندی صفحات وب به‌منظور ارتقای عملکرد موتورهای جست‌وجو مرور خواهد شد. سپس خلاصه‌وار توصیفی از الگوریتم CohsMix، که اساس مقاله‌ی حاضر است، ارائه می‌شود.

۱.۲. خوشه‌بندی صفحات وب

خروجی یک موتور جست‌وجو، علاوه بر جامعیت، باید توانایی پاسخ‌گویی مناسب به درخواست کاربر را داشته باشد. از این رو، هیرست و پدرس [۴] در سال ۱۹۹۶ ایده‌ی خوشه‌بندی صفحات وب را براساس موضوعات، برای سازمان‌دهی نتایج جست‌وجو ارائه دادند. سپس زمیر و اتزیونی (۱۹۹۸) [۲] این مسئله را مورد توجه قرار دادند. دستاورد این پژوهش‌ها، با نارسایی‌هایی همراه بود، اما توانست پنجره‌ی نوینی برای ارتقای خدمات وب به روی پژوهش‌گران بگشاید.

صرف پرداختن به موضوعات کلی که در صفحات وب مطرح می‌شود، به این دلیل که رفتارهای اجتماعی نویسندگان صفحات را در نظر نمی‌گیرد، چندان نمی‌تواند گویای گروه‌بندی واقعی و مفید برای سازمان‌دهی نتایج باشد. به عبارت دیگر، تهیه‌کنندگان محتوای صفحات مختلف اغلب در حوزه‌ی خود با افراد شبیه به خود رابطه دارند و از لینک آن‌ها در صفحات‌شان استفاده می‌کنند. بنابراین محتوایی که در یک حوزه تهیه شده می‌تواند نشأت‌گرفته از دیدگاه‌های مختلف باشد و هر دیدگاه ممکن است الگویی خاص و کاربرانی خاص داشته باشد. از این رو در سال ۲۰۰۲ هی و همکاران، [۵] بدون در نظر گرفتن محتوای صفحات، خوشه‌بندی آن‌ها را براساس لینک‌های میان صفحات انجام داده‌اند. روش دیگری نیز توسط بکرمن و همکاران (۲۰۰۶) [۶] ارائه شد که از یک روش چندعاملی^۶ برای خوشه‌بندی صفحات استفاده می‌کرد و مبنای شباهت در آن بر پایه‌ی مسیری است که میان هر دو صفحه وجود دارد.

بدین ترتیب می‌توان روش‌های بالا را در دو دسته جای داد: ۱. روش‌های مبتنی بر محتوا و موضوع صفحات؛ ۲. روش‌هایی که لینک‌های میان صفحات را عامل تشکیل گروه‌ها می‌دانند. در این میان، به نظر می‌رسد گروه‌بندی واقعی صفحات معطوف به الگوهایی است که در سایه‌ی در نظر گرفتن هم‌زمان این دو موضوع (یعنی محتوا و لینک) قابل شناسایی‌اند. چنین رویکردی توسط زنتی و همکاران (۲۰۱۰) [۱] ارائه شد؛ آنها یک الگوریتم مبتنی بر مدل شکل‌گیری گراف صفحات ارائه دادند که از روش معروف مقادیر انتظاری - بهینه‌سازی [۷] بهره می‌برد. بدین ترتیب رویکرد آن‌ها نسبت به رویکردهای دیگر (به‌عنوان مثال مدل‌های مبتنی بر فاصله) [۳،۲] توانایی بیشتری در کشف الگوهای واقعی میان صفحات داشت.

اگر مجموعه‌ی صفحات در قالب یک گراف در نظر گرفته شود، هر صفحه با

یک گره قابل نمایش است و لینک میان صفحات، یال‌های این گراف خواهد بود. زنتی و همکاران [۱] واژگان مهم و مکرر صفحات را به‌عنوان خصیصه‌های گره‌ها در نظر گرفتند که در واقع بیان‌گر محتوای صفحات است. از این منظر روش آنها را می‌توان در حوزه‌ی خوشه‌بندی گراف‌های خصیصه‌دار^۸ برشمرد. اولین پژوهش در این حوزه را نوبل و همکاران [۸] انجام دادند؛ آن‌ها تعداد خصیصه‌های با مقادیر مشابه را به‌عنوان معیار شباهت دو گره در نظر گرفتند. سپس ماتریس مجاورت وزن‌دار جدیدی با استفاده از حاصل ضرب این معیار در ماتریس مجاورت گره‌ها تعریف کردند و خوشه‌بندی براساس ماتریس جدید را به‌عنوان روشی که به‌طور هم‌زمان ویژگی‌های ساختاری و خصیصه‌های گره‌ها را لحاظ می‌کند پیشنهاد دادند. روش دیگری نیز در سال ۲۰۰۸ توسط اشتینهاسر و چاولا [۹] ارائه شد که همین رویکرد را دنبال می‌کرد. ژو و همکاران [۱۰] در سال ۲۰۰۹ یک روش مبتنی بر فاصله توسعه دادند. آن‌ها معیار جدیدی برای سنجش فاصله‌ی میان گره‌ها تعریف کردند. این معیار فاصله‌ی گره‌ها را، از نظر ساختار و خصیصه‌ها، بر مبنای فاصله‌ی قدم‌زدن تصادفی^۹ محاسبه می‌کند. در سال ۲۰۱۲ چنگ و همکاران [۱۱] این روش را با استفاده از وزن‌دهی به گره‌ها و ارائه‌ی روشی جدید برای به‌روزر کردن مرکز خوشه‌ها بهبود داده‌اند.

رویکرد مهم دیگر، رویکرد مبتنی بر مدل است. یک دسته از روش‌هایی که بر مبنای این رویکرد شکل گرفته‌اند، از مدل احتمالاتی بی‌زی^{۱۰} استفاده کرده‌اند. [۱۲،۱۳] به عبارت دیگر مسئله‌ی خوشه‌بندی گراف را به یک مسئله‌ی استنباط احتمالاتی تبدیل کرده‌اند.

۲.۲. مروری بر الگوریتم CohsMix

اساس مدل پیشنهادی در نوشتار حاضر الگوریتم موسوم به CohsMix است (شکل ۱)، که در سال ۲۰۱۰ توسط زنتی و همکاران [۱] ارائه شد و روشی مبتنی بر استنباط آماری است. الگوریتم CohsMix گونه‌ی بی‌از روش‌های تکرار

الگوریتم CohsMix	
ورودی: ماتریس مجاورت X و ماتریس خصیصه Y	
۱	تعیین مقادیر اولیه برای پارامترها $(\alpha, \pi, \sigma, \tau_{ij}, \theta)$ و σ ها به طور تصادفی
۲	در حالی که مقادیر همگرا نشده‌اند انجام بده
۳	برای هر گره $i \in \{1, 2, \dots, n\}$ انجام بده
۴	برای هر خوشه $Q \in \{1, 2, 3, \dots, Q\}$ انجام بده
۵	مقادیر τ_{ij} را با استفاده از رابطه (۱) به روز کن
۶	مقدار τ_{ij} مربوط به هر گره را بر مجموع τ_{ij} های آن گره تقسیم
۷	برای هر خوشه $Q \in \{1, 2, 3, \dots, Q\}$ انجام بده
۸	مقدار α ها را با استفاده از رابطه (۲) به روز کن
۹	برای هر خوشه $Q \in \{1, 2, 3, \dots, Q\}$ انجام بده
۱۰	مقادیر π ها را با استفاده از رابطه (۳) به روز کن
۱۱	مقادیر μ ها را با استفاده از رابطه (۴) به روز کن
۱۲	یک واحد به m اضافه کن
خروجی: مقادیر نهایی برای τ_{ij} ها و سایر پارامترها (Θ)	

شکل ۱. الگوریتم CohsMix.

حدود ۱۰ الی ۱۵ تکرار برای رسیدن به حالت بهینه‌ی محلی لازم است، محاسبات مختص هر تکرار -- به‌ویژه استفاده از رابطه‌ی به‌روzkشده‌ی τ_{iq} ها یا همان رابطه‌ی ۱ -- برای مسائل با اندازه‌ی متوسط و بزرگ زمان‌گیر است. از این رو، یافتن نقطه‌ی شروع مناسب و استفاده از روش‌های جایگزین برای کاهش محاسبات می‌تواند رویکردی مؤثر در افزایش کارایی الگوریتم CohsMix باشد. راهکار پیشنهادی در اینجا این است که خروجی الگوریتم خوشه‌بندی K-means به‌عنوان نقطه‌ی شروع نسبتاً مناسبی برای الگوریتم CohsMix تعریف شود. بدین ترتیب τ_{iq} ها و α_q های اولیه به دست می‌آید.

پیش از پرداختن به ادامه‌ی مراحل، ذکر این نکته ضروری است که زنجی و همکاران^[۱] در الگوریتم خود مقدار انحراف استاندارد توزیع نرمال مربوط به خصیصه‌های هر خوشه را برابر مقدار ثابت ۱ در نظر گرفته‌اند. چنین فرضی می‌تواند واقعی‌بودن نتایج را خدشه‌دار کند. از این رو، با توجه به پارامترها و متغیرهای تعریف شده در الگوریتم CohsMix، رابطه‌ی برای محاسبه‌ی انحراف استاندارد توزیع نرمال درباره‌ی خصیصه‌های هر خوشه به دست آورده‌ایم. رابطه‌ی ۶ رابطه‌ی مورد نظر است که در هر تکرار به روز می‌شود.

پس از یافتن مقادیر اولیه‌ی τ_{iq} ها و α_q ها براساس خروجی الگوریتم K-means، محاسبه‌ی π_{ql} ها از رابطه‌ی ۳، و محاسبه‌ی پارامترهای توزیع نرمال هر خوشه طبق روابط ۴ و ۶ ممکن خواهد بود. نتیجه‌ی استفاده از رابطه‌ی ۶، به دست آوردن یک انحراف استاندارد جداگانه برای هر خوشه خواهد بود.

$$\sigma_q^{(m+1)} = \sqrt{\frac{\sum_i \tau_{iq}^{(m+1)} (Y_i - \mu_q^{(m+1)})^T (Y_i - \mu_q^{(m+1)})}{\sum_i \tau_{iq}^{(m+1)}}}, \quad (6)$$

در پایان الگوریتم CohsMix با توجه به مقادیر جدید اجرا خواهد شد. همچنین از آنجا که عناصر ماتریس X_{ij} مقادیر صفر یا ۱ را اتخاذ می‌کند، می‌توان جمله دوم:

$$A = \prod_{j \neq i} \prod_l [(\pi_{ql}^{(m)})^{X_{ij}} (1 - \pi_{ql}^{(m)})^{1-X_{ij}}] \tau_{iq}^{(m)}, \quad (7)$$

را که حاصل ضرب عبارت‌های:

$$P = [(\pi_{ql}^{(m)})^{X_{ij}} (1 - \pi_{ql}^{(m)})^{1-X_{ij}}] \tau_{iq}^{(m)}, \quad (8)$$

است به دو بخش p_1 و p_2 تبدیل کرد که سریع‌تر قابل محاسبه‌اند. در واقع، اگر گره‌های i و j با یکدیگر رابطه داشته باشند (یعنی $X_{ij} = 1$) آنگاه:

$$P_1 = \pi_{ql}^{(m) \tau_{iq}^{(m)}}, \quad (9)$$

و در غیر این صورت:

$$P_2 = (1 - \pi_{ql}^{(m)}) \tau_{iq}^{(m)}, \quad (10)$$

استفاده از دو رابطه‌ی ۹ و ۱۰ مستلزم این است که مقادیر π_{ql} صفر یا ۱ نباشد که این فرض نیز دور از واقعیت نیست. بدین ترتیب نوآوری‌های مدل پیشنهادی عبارت‌اند از: ۱. جلوگیری از اجراهای مکرر با ایجاد جواب‌های اولیه‌ی مناسب حاصل به‌کارگیری الگوریتم K-means؛ ۲. کاهش زمان اجرای الگوریتم مورد نظر با تعریف روابط ۹ و ۱۰؛ ۳. افزایش کیفیت جواب‌ها با استفاده از محاسبه‌ی انحراف استاندارد واقعی برای هر خوشه (شکل ۲).

شونده‌ی EM^[۷] است که در آن زنجی و همکاران مفروضات مدل ساده‌ی گراف کاملاً تصادفی اردوس و رینی^[۱۲] را به‌عنوان مفروضات مدل خود پذیرفتند. به‌عبارت دیگر آن‌ها فرض کردند که یال‌های یک شبکه به‌طور تصادفی و با احتمال مشخص ایجاد شده است. مشخصه‌ی دیگر این الگوریتم این است که توزیع احتمالی خصیصه‌های گره‌های هم‌خوشه را نرمال در نظر می‌گیرد.

در مورد گراف یا شبکه G فرض کنید که X_{ij} نمایان‌گر درایه‌های ماتریس مجاورت گره‌ها و Y_{ik} درایه‌های ماتریس مربوط به مقادیر K ویژگی (ویژگی $k = 1, 2, \dots, K$) هر یک از گره‌ها باشد. الگوریتم CohsMix (شکل ۱) با در نظر گرفتن ماتریس‌های X و Y تلاش دارد همه‌ی N گره موجود را به Q خوشه تخصیص دهد. به‌عبارت دیگر این روش از توزیع احتمالاتی X و Y برای تخمین میزان وابستگی هر گره به هر یک از خوشه استفاده می‌کند. این میزان وابستگی با τ_{iq} نشان داده می‌شود که در آن $Q = 1, 2, 3, \dots, Q$. در ابتدا مقادیر اولیه‌ی τ_{iq} ها به‌طور تصادفی تعیین می‌شود اما در ادامه، مقادیر به‌دست آمده در هر مرحله از الگوریتم روش CohsMix توسط رابطه‌ی ۱ به‌روز می‌شود:

$$\tau_{iq}^{(m+1)} = \alpha_q^{(m)} \prod_{j \neq i} \prod_l [(\pi_{ql}^{(m)})^{X_{ij}} (1 - \pi_{ql}^{(m)})^{1-X_{ij}}] \tau_{iq}^{(m)} \times \prod_{k=1}^K [\exp(-\frac{1}{2\sigma_q^{(m)}} (Y_{ik} - \mu_{qk}^{(m)})^2)] \quad (1)$$

که در آن:

$$\alpha_q^{(m+1)} = \frac{1}{N} \sum_{i=1}^N \tau_{iq}^{(m+1)}, \quad (2)$$

$$\pi_{ql}^{(m+1)} = \frac{\sum_{i \neq j} \tau_{iq}^{(m+1)} \tau_{jl}^{(m+1)} X_{ij}}{\sum_{i \neq j} \tau_{iq}^{(m+1)} \tau_{jl}^{(m+1)}}, \quad (3)$$

$$\mu_q^{(m+1)} = \frac{\sum_i \tau_{iq}^{(m+1)} Y_i}{\sum_i \tau_{iq}^{(m+1)}}, \quad (4)$$

پیچیدگی زمانی این الگوریتم از مرتبه‌ی $O(N^2)$ است. تعداد بهینه‌ی خوشه‌ها (یا همان Q) با اجرای الگوریتم برای مقادیر متفاوت Q و با توجه به بیشینه‌ی مقدار به دست آمده برای هر Q با استفاده از معیار ICL^[۱۵] به دست می‌آید:

$$ICL(Q) = \max_{\Theta} \log L(X, Y, Z; \Theta, Q) - \frac{1}{4} Q(Q-1) \log\left(\frac{N(N-1)}{4}\right) - \frac{Q-1}{4} \log(N) - K(K-1) \log\left(\frac{N(N-1)}{4}\right) + KQ \log\left(\frac{N(N-1)}{4}\right) \quad (5)$$

جمله‌ی اول معیار ICL بیشینه‌ی لگاریتم راستی‌نمایی ماتریس‌های ورودی، و Z متغیر صفر و ۱ است که نمایان‌گر تخصیص یا عدم تخصیص یک گره به هر یک از Q خوشه است.

۳. الگوریتم پیشنهادی

الگوریتم CohsMix با تعداد اندکی تکرار به حالت بهینه‌ی محلی خود می‌رسد و سازوکاری برای فرار از بهینه‌ی محلی ندارد. بنابراین برای رسیدن به یک جواب قابل قبول لازم است چندین بار الگوریتم اجرا شود و بهترین اجرا ثبت شود. اگرچه فقط

جدول ۱. اطلاعات مربوط به پارامترهای مجموعه داده‌ها.

پارامترها				مسائل	پارامترها				مسائل
σ	$d(\mu_q, \mu_i)$	Q	N		σ	$d(\mu_q, \mu_i)$	Q	N	
۲	۱	۴	۱۵۰	۱۱	۲	۳	۳	۱۰۰	۱
۲	۲	۴	۱۵۰	۱۲	۲	۳	۴	۱۰۰	۲
۲	۳	۴	۱۵۰	۱۳	۲	۳	۵	۱۰۰	۳
۲	۴	۴	۱۵۰	۱۴	۲	۳	۶	۱۰۰	۴
۲	۵	۴	۱۵۰	۱۵	۲	۳	۷	۱۰۰	۵
۱٫۵	۵	۴	۱۵۰	۱۶	۲	۳	۳	۱۵۰	۶
۲٫۵	۵	۴	۱۵۰	۱۷	۲	۳	۳	۲۰۰	۷
۳	۵	۴	۱۵۰	۱۸	۲	۳	۳	۲۵۰	۸
۳٫۵	۵	۴	۱۵۰	۱۹	۲	۳	۳	۳۵۰	۹
					۲	۳	۳	۵۰۰	۱۰

جدول ۲. نتایج مربوط به زمان اجرای دو الگوریتم روی مسائل با تعداد گره متفاوت برحسب ثانیه.

روش‌ها	مسائل					
	۱۰	۹	۸	۷	۶	۱
CohsMix	۵۰٫۵۹	۲۵٫۰۸	۱۲٫۹	۸٫۴	۴٫۷۴	۲٫۱۲
الگوریتم جدید	۳۴٫۸۹	۱۷٫۸۱	۹٫۱۹	۵٫۸۵	۳٫۳۲	۱٫۵۳
اختلاف نسبی	۰٫۳۱	۰٫۲۹	۰٫۲۹	۰٫۳	۰٫۳	۰٫۲۸

برای اعداد بزرگ، توزیع نرمال در نظر گرفته شده است. در این جدول، N تعداد گره‌ها، Q بیانگر تعداد خوشه‌های از پیش تعیین شده برای تولید ماتریس‌ها و $d(\mu_i, \mu_q)$ تفاضل میان میانگین خصیصه‌های دو خوشه است. معیار فاصله برای در نظر گرفتن این تفاضل فاصله‌ی اقلیدسی لحاظ شده است. ستون σ اشاره به مقدار انحراف استاندارد میان مقادیر خصیصه‌ها دارد. همچنین پارامتر دیگری در ایجاد مجموعه داده‌ها کاربرد داشته که تفاضل میان: «احتمال برقراری ارتباط میان یک گره و گره‌های هم‌خوشه» با «احتمال برقراری ارتباط میان آن گره و گره‌های ناهم‌خوشه» را نشان می‌دهد. از آنجا که این پارامتر برای همه‌ی مجموعه داده‌ها یکسان و برابر $۰٫۲$ است در جدول مورد اشاره قرار نگرفته است. خصیصه‌ها مقادیری گویا دارند و تعداد آن‌ها برای هر گره در همه‌ی مجموعه داده‌ها برابر ۳ است. در جدول ۱ اطلاعات مربوط به مجموعه داده‌های مختلف ارائه شده است.

۲.۱.۴. مقایسه‌ی روش‌ها از نظر زمان اجرا

در این بخش ابتدا دو الگوریتم (الگوریتم CohsMix و الگوریتم جدید پیشنهادی) از نظر زمان اجرا برحسب ثانیه روی مسائلی که تعداد خوشه‌ی متفاوت دارند مقایسه می‌شود. انتظار می‌رود نوآوری‌های مربوط به کاهش حجم محاسبات بتواند بهبود قابل توجهی به هنگام رشد مسئله نشان دهد. در ادامه، زمان اجرای دو الگوریتم روی مسائل متفاوت از نظر تعداد خوشه $d(\mu_q, \mu_i)$ و انحراف استاندارد با هم مقایسه خواهد شد. لازم به یادآوری است که نتایج ارائه شده برای الگوریتم CohsMix میانگین نتایج حاصل از ۱۰ بار اجرای این الگوریتم است.

نتایج مربوط به استفاده از دو روش در جدول ۲ آمده است. طبق این جدول تفاوت قابل ملاحظه‌ی میان زمان اجرای دو روش رؤیت می‌شود. مشخص است که با افزایش تعداد گره‌ها (از ۱۰۰ به ۵۰۰ گره طبق جدول ۱)، این اختلاف نیز بیشتر می‌شود. از سوی دیگر می‌توان نتیجه گرفت که الگوریتم جدید همواره حدود

الگوریتم جدید	
ورودی: ماتریس مجاورت X و ماتریس خصیصه‌ها Y	
۱ خوشه بندی ماتریس خصیصه‌ها با استفاده از الگوریتم K-means	
۲ محاسبه مقادیر $(\alpha, \pi, \tau_{iq}, \sigma)$ ها با توجه به خروجی گام قبل	
۳ در حالی که مقادیر همگرا نشده اند انجام بده	
۴ برای هر گره $i \in \{1, 2, \dots, N\}$ انجام بده	
۵ برای هر خوشه $Q \in \{1, 2, \dots, Q\}$ انجام بده	
۶ مقادیر τ_{iq} را با استفاده از روابط (۱)، (۹) و (۱۰) به روز کن	
۷ مقدار τ_{iq} مربوط به هر گره را بر مجموع τ_{iq} های آن گره تقسیم کن	
۸ برای هر خوشه $Q \in \{1, 2, \dots, Q\}$ انجام بده	
۹ مقدار α ها را با استفاده از رابطه (۲) به روز کن	
۱۰ برای هر خوشه $I \in \{1, 2, \dots, Q\}$ انجام بده	
۱۱ مقادیر π ها را با استفاده از رابطه (۳) به روز کن	
۱۲ مقادیر μ ها را با استفاده از رابطه (۴) به روز کن	
۱۳ مقادیر μ ها را با استفاده از رابطه (۶) به روز کن	
۱۴ یک واحد به m اضافه کن	
خروجی: مقادیر نهایی برای τ_{iq} ها و سایر پارامترها (Θ)	

شکل ۲. الگوریتم جدید.

۴. نتایج عددی

در این بخش ابتدا داده‌های مورد استفاده معرفی می‌شوند. سپس نتایج اجرای دو روش (روش اصلی CohsMix و روش پیشنهادی) روی مجموعه داده‌های مختلف ارائه و تأثیر استفاده از نوآوری‌های روش پیشنهادی برای کاهش زمان اجرای الگوریتم و افزایش کیفیت جواب‌ها بررسی می‌شود. لازم به ذکر است از آنجا که الگوریتم CohsMix با در نظر گرفتن نقاط شروع تصادفی اجرا می‌شود، برای این که مقایسه‌ی منصفانه‌ی میان نتایج دو الگوریتم داشته باشیم، لازم است الگوریتم CohsMix (الگوریتم پایه) بیشتر از یک بار اجرا شود و میانگین نتایج به دست آمده از ۱۰ بار اجرای این الگوریتم در مقایسه و ارزیابی در نظر گرفته شود.

۱.۴. آزمایش روی مجموعه داده‌های مصنوعی

یکی از روش‌های مقایسه‌ی الگوریتم‌ها استفاده از مجموعه داده‌ی است که اطلاعات کافی نسبت به خوشه‌های موجود در آن دسترس باشد. با داشتن چنین داده‌هایی می‌توان تعیین کرد که خروجی کدام یک از الگوریتم‌ها بیشتر شبیه آن چیزی است که در مجموعه داده وجود دارد. بدین منظور در بخش بعد چنین مجموعه داده‌ی را معرفی می‌کنیم و سپس به بررسی نتایج خواهیم پرداخت.

۱.۱.۴. مجموعه داده‌های مصنوعی

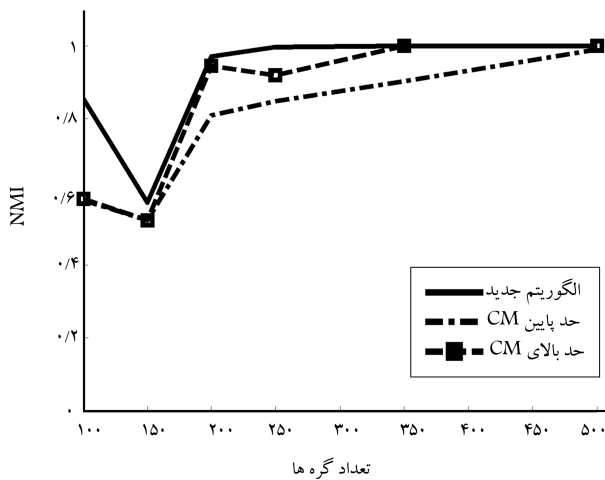
مجموعه داده‌های مورد استفاده در این نوشتار از مطالعه‌ی زنقی و همکاران^[۱] انتخاب شده است. این مجموعه داده‌ها، ماتریس‌های شبیه‌سازی شده توسط رایانه است که با در نظر گرفتن پارامترهای زیر ایجاد شده‌اند. شبیه‌سازی ماتریس مجاورت X برای این مجموعه داده‌ها با به‌کارگیری فرایند ایجاد گراف تصادفی و نیز توزیع آماری خصیصه‌ها، با استناد به قضیه‌ی حد مرکزی

شکل کمی بیان کند. برای این مهم، معیار NMI^[۱۶] در نظر گرفته می‌شود. محاسبه‌ی این معیار برای دو مجموعه‌ی A و B مطابق رابطه‌ی ۱۱ انجام می‌شود:

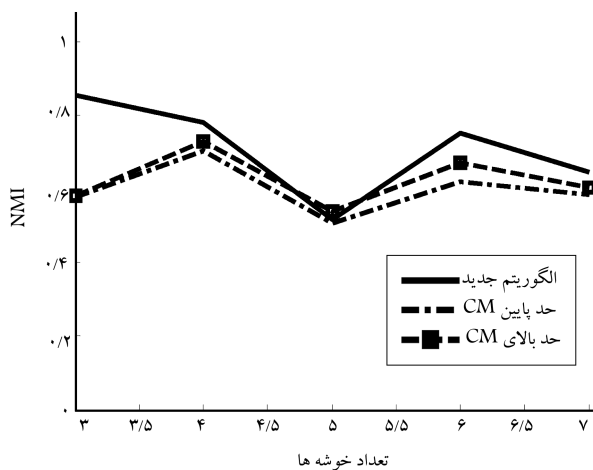
$$I(A, B) = \frac{-2 \sum_{i=1}^{C_A} \sum_{j=1}^{C_B} \frac{C_{ij} \log(C_{ij} N)}{C_i \cdot C_j}}{\sum_{i=1}^{C_A} C_i \cdot \log\left(\frac{C_i}{N}\right) + \sum_{j=1}^{C_B} C_j \cdot \log\left(\frac{C_j}{N}\right)}, \quad (11)$$

که در آن C_A و C_B به ترتیب نماینده‌ی تعداد خوشه‌ها در دو مجموعه‌ی A و B هستند و C_i ، C_j و C_{ij} نیز به ترتیب بیان‌گر تعداد گره‌ها در خوشه‌ی i ام، j ام و مشترک بین i ام و j ام است. هرچه میزان NMI بیشتر باشد خوشه‌بندی با کیفیت بالاتری حاصل شده است. در مسائل مورد بحث در نوشتار حاضر استفاده از معیار NMI امکان‌پذیر خواهد بود، زیرا مجموعه داده‌ها توسط رایانه ایجاد شده و اطلاعاتی در این خصوص که کدام‌یک از گره‌ها با هم هم‌خوشه‌اند وجود دارد.

نتایج به دست آمده از اجرای دو الگوریتم روی مجموعه داده‌های مختلف برحسب معیار NMI در شکل‌های ۳ تا ۶ نشان داده شده است. در این شکل‌ها، الگوریتم جدید با کم‌ترین و بیشترین مقدار به دست آمده از ۱۰ اجرا الگوریتم CohsMix (که به ترتیب با عنوان حد پایین CM و حد بالا CM نشان داده شده‌اند) مقایسه می‌شود.



شکل ۳. نتیجه‌ی اجرای دو الگوریتم روی مسائل با تعداد گره متفاوت براساس معیار NMI.



شکل ۴. نتیجه‌ی اجرای دو الگوریتم روی مسائل با تعداد خوشه‌ی متفاوت براساس معیار NMI.

جدول ۳. نتایج مربوط به زمان اجرای دو الگوریتم روی مسائل با تعداد خوشه متفاوت برحسب ثانیه.

روش‌ها	مسائل				
	۵	۴	۳	۲	۱
CohsMix	۹٫۱۱	۶٫۸۸	۵٫۰۵	۳٫۴۴	۲٫۱۲
الگوریتم جدید	۶٫۲۴	۴٫۶۶	۳٫۴۷	۲٫۴	۱٫۵۳
اختلاف نسبی	۰٫۳۲	۰٫۳۲	۰٫۳۱	۰٫۳	۰٫۲۸

جدول ۴. نتایج مربوط به زمان اجرای دو الگوریتم روی مسائل با $d(\mu_q, \mu_l)$ متفاوت برحسب ثانیه.

روش‌ها	مسائل				
	۱۵	۱۴	۱۳	۱۲	۱۱
CohsMix	۷٫۶۳	۷٫۶۵	۷٫۷۹	۷٫۷۲	۷٫۷۴
الگوریتم جدید	۵٫۳	۵٫۳	۵٫۳۱	۵٫۳۳	۵٫۲۹
اختلاف نسبی	۰٫۳۱	۰٫۳۱	۰٫۳۲	۰٫۳۱	۰٫۳۲

جدول ۵. نتایج اجرای دو الگوریتم روی مسائل با انحراف استاندارد متفاوت برحسب ثانیه.

روش‌ها	مسائل				
	۱۹	۱۸	۱۷	۱۵	۱۶
CohsMix	۷٫۵	۷٫۵۷	۷٫۵۲	۷٫۶۳	۷٫۴۵
الگوریتم جدید	۵٫۲۶	۵٫۳۸	۵٫۳۲	۵٫۳	۵٫۳
اختلاف نسبی	۰٫۳	۰٫۲۹	۰٫۲۹	۰٫۳۱	۰٫۲۹

۳۰ درصد زمان اجرا را بهبود داده است. چنان‌که اشاره شد، این بهبود در وهله‌ی اول معطوف به کاهش تعداد تکرارهاست که در نتیجه‌ی یافتن نقطه‌ی شروع مناسب با استفاده از الگوریتم خوشه‌بندی صورت می‌پذیرد. در وهله‌ی دوم کاهش زمان انجام هر تکرار الگوریتم که در نتیجه‌ی ساده‌سازی محاسبات مرتبط با به‌روز کردن متغیرها حاصل می‌شود عاملی برای کاهش زمان حل خواهد بود. بدین ترتیب تأثیر این تغییرات در مجموعه داده‌های مختلف تقریباً مشابه بوده و از نظر زمانی قابل توجه است.

در جدول ۳ زمان اجرای دو الگوریتم، هنگامی که تعداد خوشه‌های تعریف شده در مجموعه داده‌های مصنوعی متفاوت است، ارائه شده است. این جدول نیز بیانگر کاهش زمان اجرا به‌هنگام استفاده از الگوریتم جدید است. همچنین اختلاف نسبی زمان اجرای دو الگوریتم با افزایش تعداد خوشه‌ها (از ۳ به ۷ خوشه طبق جدول ۱) رشد کمی را نشان می‌دهد.

زمان اجرای دو الگوریتم روی مجموعه داده‌های متفاوت از نظر $d(\mu_q, \mu_l)$ و انحراف استاندارد به ترتیب در جدول‌های ۴ و ۵ ارائه شده است. مشخص است که تفاوت این دو پارامتر، از آنجا که در پیچیدگی زمانی الگوریتم تأثیری ندارند، تفاوت محسوسی در زمان اجرای دو الگوریتم نخواهد داشت. از طرف دیگر، نتایج به دست آمده از زمان اجرای الگوریتم جدید در هر دو جدول کم‌تر از روش CohsMix گزارش شده است.

۳.۱.۴. مقایسه‌ی روش‌ها از نظر کیفیت خوشه‌بندی

برای مقایسه‌ی کیفیت پاسخ‌های دو روش نیاز به معیاری است که با توجه به مشخص بودن برحسب هر گره، میزان صحت خوشه‌بندی انجام گرفته توسط هر روش را به

البته به نظر می‌رسد زمانی که تفاضل میان میانگین مقادیر مربوط به خصیصه‌ها افزایش می‌یابد، هر دو الگوریتم میل به کشف صحیح خوشه‌ها دارند. علت آن است که افزایش فاصله‌ی میان مقادیر مربوط به خصیصه‌ها به معنای افزایش میزان جدایی خوشه‌ها از هم و افزایش توانایی تشخیص از یکدیگر خواهد بود.

در شکل ۶ دو الگوریتم مبتنی بر معیار NMI روی مسائل با مقدار انحراف استاندارد متفاوت مقایسه شده است. لازم به یادآوری است که روش CohsMix همواره مقدار انحراف استاندارد همه‌ی خصیصه‌ها برای همه‌ی خوشه‌ها را برابر ۱ در نظر می‌گیرد. از این رو انتظار می‌رود که با افزایش انحراف استاندارد واقعی خصیصه‌ها، کیفیت خوشه‌بندی در این روش کاهش یابد (شکل ۶). از طرفی انتظار می‌رود الگوریتم جدید که همواره برای همه‌ی خوشه‌ها مقدار انحراف استاندارد را محاسبه می‌کند، کم‌تر متأثر از تغییر انحراف استاندارد باشد. در شکل ۶ فقط زمانی که انحراف استاندارد برابر ۳ شده است میزان NMI محاسبه شده برای الگوریتم جدید کم‌تر از ۰٫۹ است. بدین ترتیب می‌توان نتیجه گرفت که نوآوری این مقاله در راستای استفاده از رابطه‌ی ۶ برای محاسبه‌ی انحراف استاندارد توزیع نرمال مربوط به خصیصه‌های خوشه‌ها مؤثر بوده است.

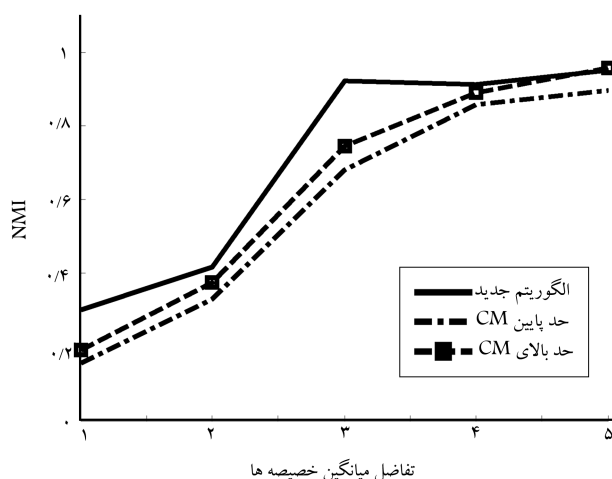
۲.۴. مطالعه‌ی موردی: وبلاگ‌های فارسی

براساس رتبه‌بندی سایت الکسا (Alexa.com) در بهار سال ۲۰۱۴ میلادی، ۳ مورد از ۱۰ سایت پر بازدید توسط کاربران ایرانی مربوط به سایت‌های ارائه‌دهنده‌ی خدمات وبلاگ‌نویسی است. این موضوع به‌خوبی جایگاه وبلاگ‌نویسی در میان کاربران ایرانی و اهمیت مطالعه روی این حوزه را نمایان می‌کند. در این بخش ابتدا به تبیین نحوه‌ی تهیه‌ی داده‌ها، و سپس به بررسی نتایج حاصل از اجرای الگوریتم پیشنهادی روی داده‌های جمع‌آوری شده خواهیم پرداخت.

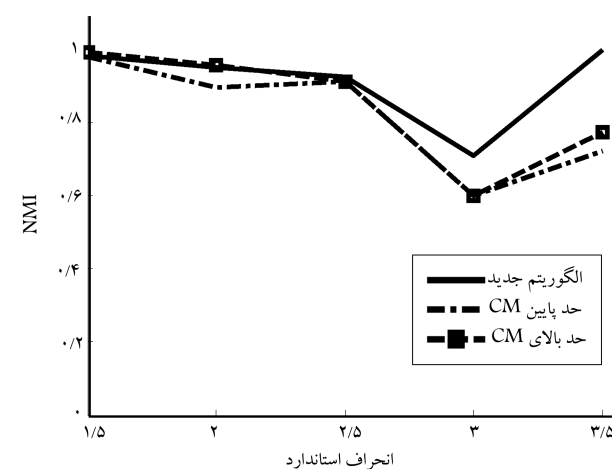
۲.۴.۱. تهیه‌ی مجموعه داده

برای تهیه‌ی داده‌های مذکور لازم است هر دو موضوع ارتباط میان وبلاگ‌ها و محتوای آن‌ها در نظر گرفته شود. به عبارت دیگر ما با شبکه‌ی از وبلاگ‌ها مواجهیم و برای جمع‌آوری داده‌ها باید از روش‌های جمع‌آوری داده که در حوزه‌ی شبکه‌های پیچیده مطرح است استفاده شود. از این رو به نظر می‌رسد استفاده از روش خزش در شبکه مناسب باشد. بدین ترتیب ابتدا به‌طور تصادفی یک وبلاگ را از یکی از سایت‌های ارائه‌دهنده‌ی خدمات وبلاگ‌نویسی برمی‌گزینیم. در اینجا سایت پرشین بلاگ (Persianblog.ir) مورد توجه واقع شده است. سپس وبلاگ‌هایی که در لیست لینک‌های این وبلاگ (دوستان نویسنده‌ی وبلاگ) هستند انتخاب می‌شوند. این نحوه‌ی گزینش در باره‌ی دوستان نویسندگان وبلاگ‌هایی که انتخاب شده‌اند نیز صورت می‌گیرد. لازم به ذکر است که وبلاگ دوستان بعضی از وبلاگ‌های مورد مطالعه در دامنه‌ی سایت‌های دیگر تعریف شده است. در اینجا از در نظر گرفتن چنین وبلاگ‌هایی صرف نظر شده زیرا هدف مطالعه و بررسی وبلاگ‌های تعریف شده در دامنه‌ی پرشین بلاگ است.

با لحاظ کردن موارد بالا تعداد ۱۶۲ وبلاگ در نهایت گزینش شد که ۳۹۰ لینک با یکدیگر دارند. محتوای این وبلاگ‌ها، متن پست‌هایی است که در وبلاگ قرار داده شده است. با استفاده از ابزار متن‌کاوی، واژگان به کار رفته و واژگان مهم متون وبلاگ‌ها استخراج شد. تنوع کل واژگان ۱۲۹۰۶ واژه بوده که براساس روش متن‌کاوی ارائه شده^[۱۷] تعداد ۸۸ مورد حائز اهمیت تشخیص داده شد. بدین ترتیب، خصیصه‌های صفحات در واقع همان واژگان مهم و مقدار آن‌ها همان فراوانی آن‌ها در صفحه‌ی مورد نظر خواهد بود.



شکل ۵. نتیجه‌ی اجرای الگوریتم‌ها براساس NMI روی مسائل با مقدار متفاوت تفاضل میانگین خصیصه‌ها.



شکل ۶. نتیجه‌ی اجرای دو الگوریتم روی مسائل با مقدار انحراف استاندارد متفاوت براساس معیار NMI.

در شکل ۳ نتیجه‌ی اجرای دو الگوریتم روی مسائل با تعداد گره متفاوت براساس معیار NMI نشان داده شده است. چنان که مشاهده می‌شود در اغلب مسائل، الگوریتم جدید خوشه‌بندی با کیفیت بالاتری را به ارمغان آورده است. از سوی دیگر، با افزایش تعداد گره‌ها و در نتیجه افزایش تعداد اعضای هر خوشه، تقریب‌های بهتری از توزیع نرمال خصیصه‌های اعضای هر خوشه حاصل شده و در نتیجه کیفیت خوشه‌بندی برحسب معیار NMI ارتقاء یافته است.

در شکل ۴ رفتار دو الگوریتم هنگام اجرا روی مجموعه داده‌های متفاوت از نظر تعداد خوشه نشان داده شده است. هر دو الگوریتم بر اثر افزایش تعداد خوشه (بدون در نظر گرفتن مقادیر متناظر با تعداد خوشه برابر ۵) سیری نزولی از نظر میزان NMI محاسبه شده دارند. به‌طور کلی انتظار می‌رود برای مسائل با تعداد خوشه‌ی بالاتر از ۴، میزان کیفیت خوشه‌بندی‌های دو الگوریتم از ۰٫۸ کم‌تر باشد. از سوی دیگر، چنان که در این شکل مشاهده می‌شود، الگوریتم جدید از نظر کیفیت جواب‌ها در اکثر مسائل برتر است.

در شکل ۵ روندی صعودی برای هر دو الگوریتم هنگام اجرا روی مجموعه داده‌های متفاوت از نظر تفاضل میانگین خصیصه‌های خوشه‌ها نشان داده شده است. مشخص است که الگوریتم جدید از نظر کیفیت جواب‌ها در اکثر مسائل برتر از روش پایه است.

کاربران پرشین بلاگ است. اگر مسئله، یافتن صفحه‌ی مرتبط با پرسش کاربر از میان صفحات موجود (مورد بررسی در داده‌های جمع‌آوری شده) باشد آنگاه ابتدا باید همه‌ی خوشه‌هایی که صفحه‌ی مرتبط با آن پرسش دارند مشخص شود. سپس از هریک از آن‌ها گزینه‌ی انتخاب می‌شود. بدین ترتیب در خروجی جست‌وجو، انواع صفحات مرتبط لحاظ شده است.

۵. نتیجه‌گیری

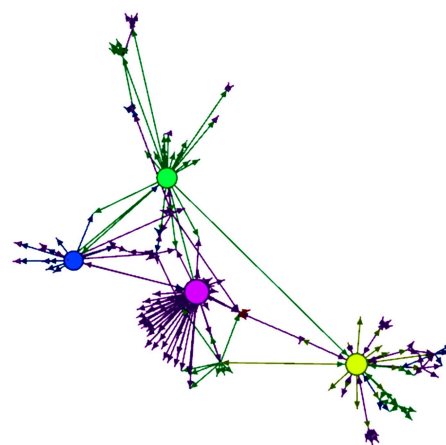
در این نوشتار مسئله‌ی خوشه‌بندی صفحات وب براساس محتوا و لینک بررسی شد. برای حل چنین مسئله‌ی یکی از الگوریتم‌های موجود در ادبیات به کار گرفته شد و تلاش شد تا با اتخاذ نوآوری‌های پیشنهادی، بهبود قابل ملاحظه‌ی -- چه از لحاظ کیفیت پاسخ‌ها و چه از لحاظ زمان حل -- ایجاد شود. این نوآوری‌ها در قالب تعیین جواب اولیه‌ی مناسب با به‌کارگیری روش خوشه‌بندی K-means، کاهش حجم محاسبات مربوط به پارامترها با استفاده از ویژگی‌های شبکه‌ی لینک‌ها و به‌کارگیری مقادیر واقعی انحراف استانداردهای هر خوشه بیان شد. نتایج به دست آمده از انجام آزمایش‌های متنوع روی مجموعه داده‌های استاندارد حاکی از عملکرد مناسب روش پیشنهادی نسبت به روش CohsMix است. علاوه بر این، مجموعه داده‌ی مربوط به وبلاگ‌های فارسی استخراج شد و الگوریتم پیشنهادی بر روی آن داده به‌کار گرفته شد.

پانویس‌ها

1. covariates on hidden structure using mixture models (CohsMix)
2. query
3. clustering
4. content
5. link
6. multi-agent
7. expectation-maximization
8. attributed graphs
9. random walk
10. bayesian probabilistic model
11. integrated classification likelihood (ICL)
12. normalized mutual index

منابع (References)

1. Zanghi, H., Volant, S. and Ambroise, C. "Clustering based on random graph model embedding vertex features", *Pattern Recogn. Lett.*, **31**(9), pp. 830-836 (2010).
2. Weiss, R., Velez, B. and Sheldon, M. "HyPursuit: A hierarchical network search engine that exploits content-link hypertext clustering", *Hypertext'96 Proceedings of the Seventh ACM Conference on Hypertext*, New York, USA, pp. 180-193 (1996).
3. Zamir, O. and Etzioni, O. "Web document clustering: A feasibility demonstration", In: *Proc. 21st Annual Inter-*



شکل ۷. نمایش خوشه‌های مربوط به وبلاگ‌ها.

۲.۲.۴. نتایج مطالعه‌ی موردی

تعداد بهینه‌ی خوشه‌ها براساس خروجی الگوریتم جدید برابر ۴ به دست آمد. در شکل ۷ تعداد ۴ خوشه کشف شده حول مراکز خوشه و نیز صفحاتی که مقادیر به دست آمده برای خصیصه‌های آن‌ها به میانگین مقادیر صفحات هم‌خوشه‌ی آن‌ها شبیه‌ترین است با دوائر بزرگ‌تر نشان داده شده‌اند. به‌طور کلی این شکل بیان‌گر تنوع مطالب در میان صفحات جمع‌آوری شده از

nat. ACM SIGIR Conf. on Research and Development in Information Retrieval, New York, USA, pp. 46-54 (1998).

4. Hearst, M. and Pedersen, J. "Reexamining the cluster hypothesis: Scatter/gather on retrieval results", *19th Annual Internet. ACM SIGIR Conf. Research and Development in Information Retrieval*, New York, USA, pp. 76-84 (1996).
5. He, X., Zha, H., Ding, C.H.Q. and Simon, H.D. "Web document clustering using hyperlink structures", *Computat. Statist. Data Anal.*, **41**(1), pp. 19-45 (2002).
6. Bekkerman, R., Ziberstein, S. and Allan, J. "Web page clustering using heuristic search in the web graph", *IJCAI'07 Proceedings of the 20th International Joint Conference on Artificial Intelligence*, New York, USA, pp. 2280-2285 (2006).
7. Dempster, A.P., Laird, N.M. and Rubin, D.B. "Maximum likelihood from incomplete data via the EM algorithm", *J. R. Stat. Soc.*, **39**(1), pp. 1-38 (1977).
8. Neville, J., Adler, M. and Jensen, D. "Clustering relational data using attribute and link information", In *Proceedings of the Text Mining and Link Analysis Workshop, 18th International Joint Conference on Artificial Intelligence*, pp. 689-698 (2003).
9. Steinhaeuser, K. and Chawla, V. "Community detection in a large real-world social network", *Soc. Comput. Beh. Modeling*, pp. 168-175 (2008).

10. Zhou, Y., Cheng, H. and Yu, J.X. "Graph clustering based on structural/attribute similarities", *Proceedings of the VLDB Endowment*, pp. 718-729 (2009).
11. Cheng, H., Zhou, Y., Huang, X. and Yu, J.X. "Clustering large attributed information networks: An efficient incremental computing approach", *Data Min. Knowl. Disc.*, **25**, pp. 450-477 (2012).
12. Kim, M. and Leskovec, J. "Multiplicative attribute graph model of real-world networks", *Internet Math.*, **8**, pp. 113-160 (2012).
13. Xu, Z., Ke, Y., Wang, Y., Cheng, H. and Cheng, J. "A model-based approach to attributed graph clustering", *SIGMOD'12 International Conference on Management of Data*, Scottsdale, Arizona, USA, pp. 505-516 (2012).
14. Erdős, P. and Rényi, A. "On random graphs", *Publicationes Mathematicae*, **6**, pp. 290-297 (1959).
15. Biernacki, C., Celeux, G. and Govaert, G. "Assessing a mixture model for clustering with the integrated completed likelihood", *IEEE PAMI*, **22**(7), pp. 719-725 (2000).
16. Danon, L., Diaz-Guilera, A., Duch, J. and Arenas, A. "Comparing community structure identification", *Journal of Statistical Mechanics: Theory and Experiment*, **09**, pp. 1-10 (2005)
17. Banchs, R., *Text Mining with MATLAB*, 1st Edn., chapter 8, Springer (2013).